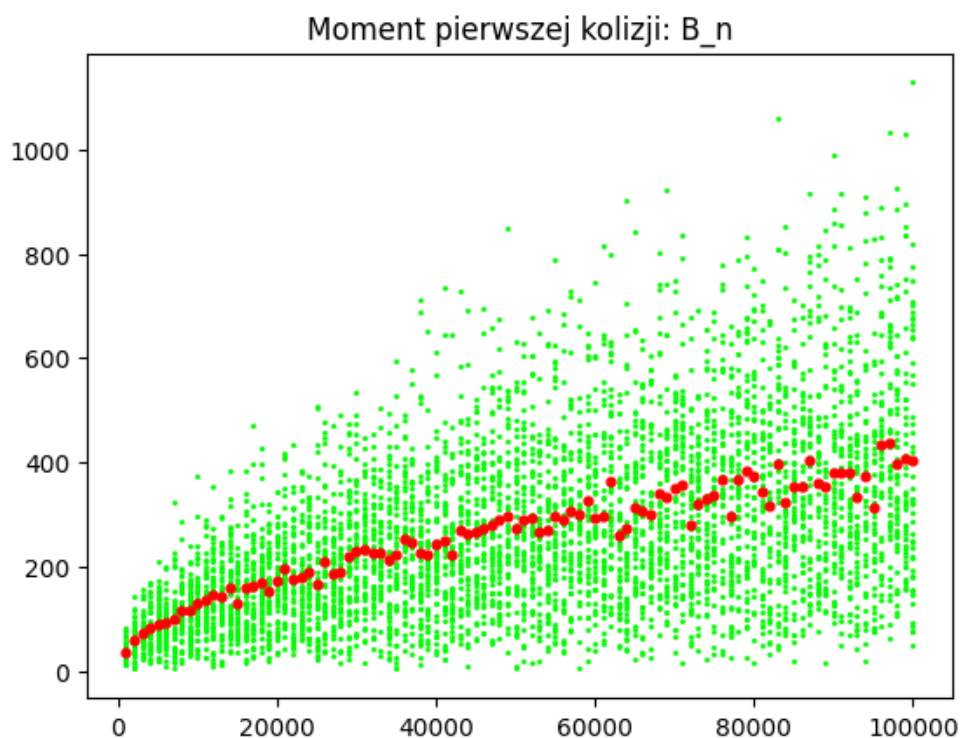
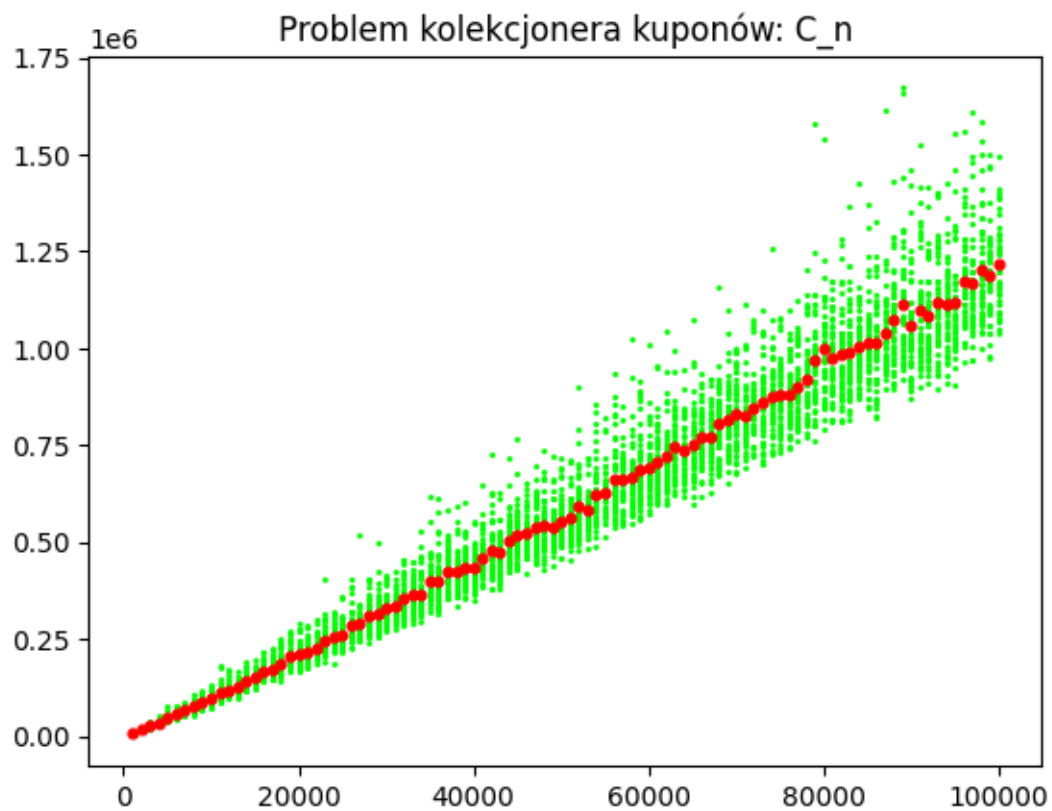
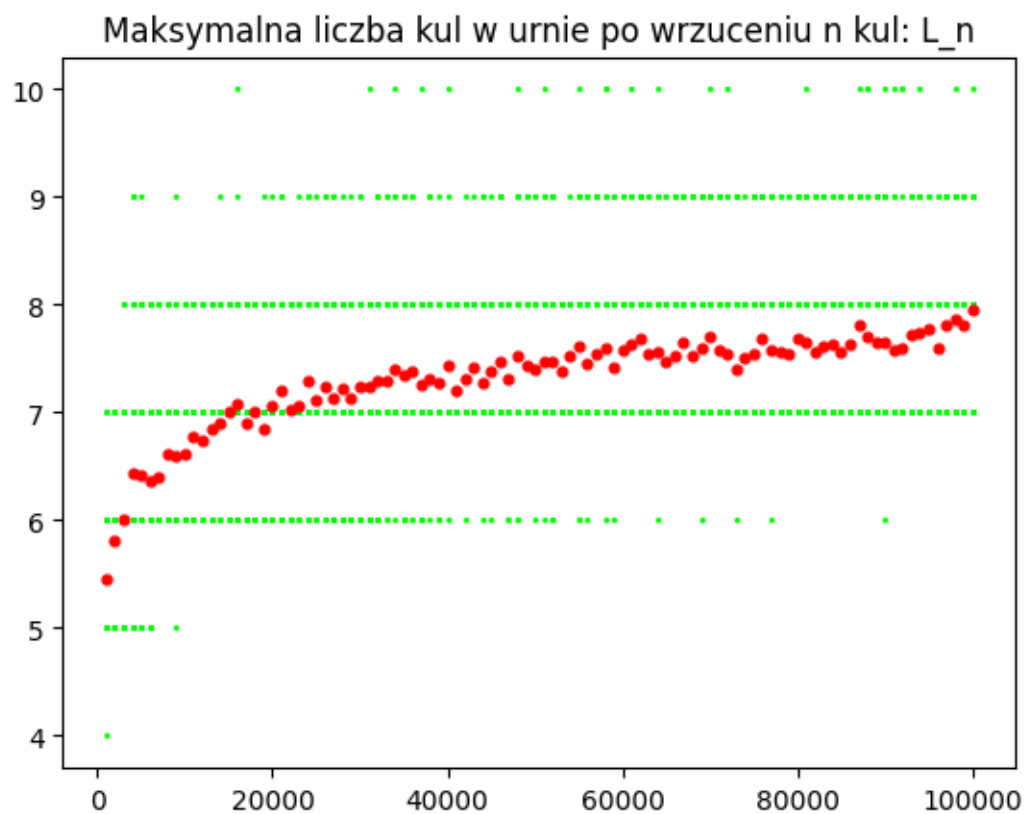
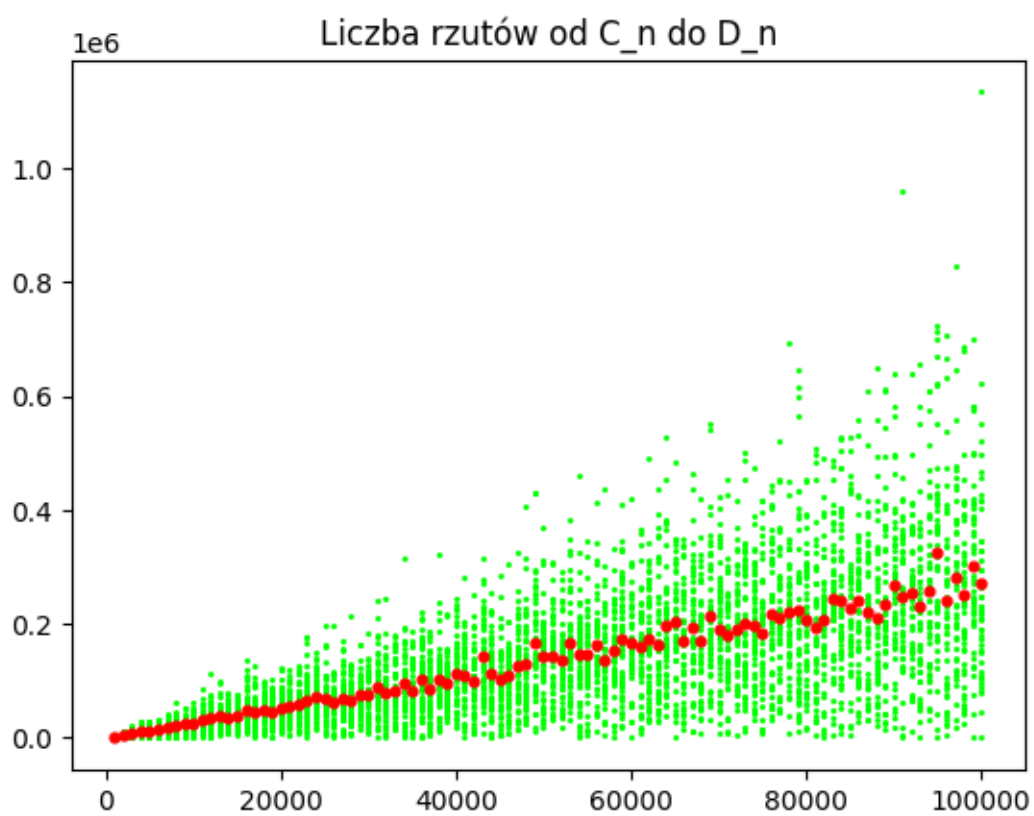
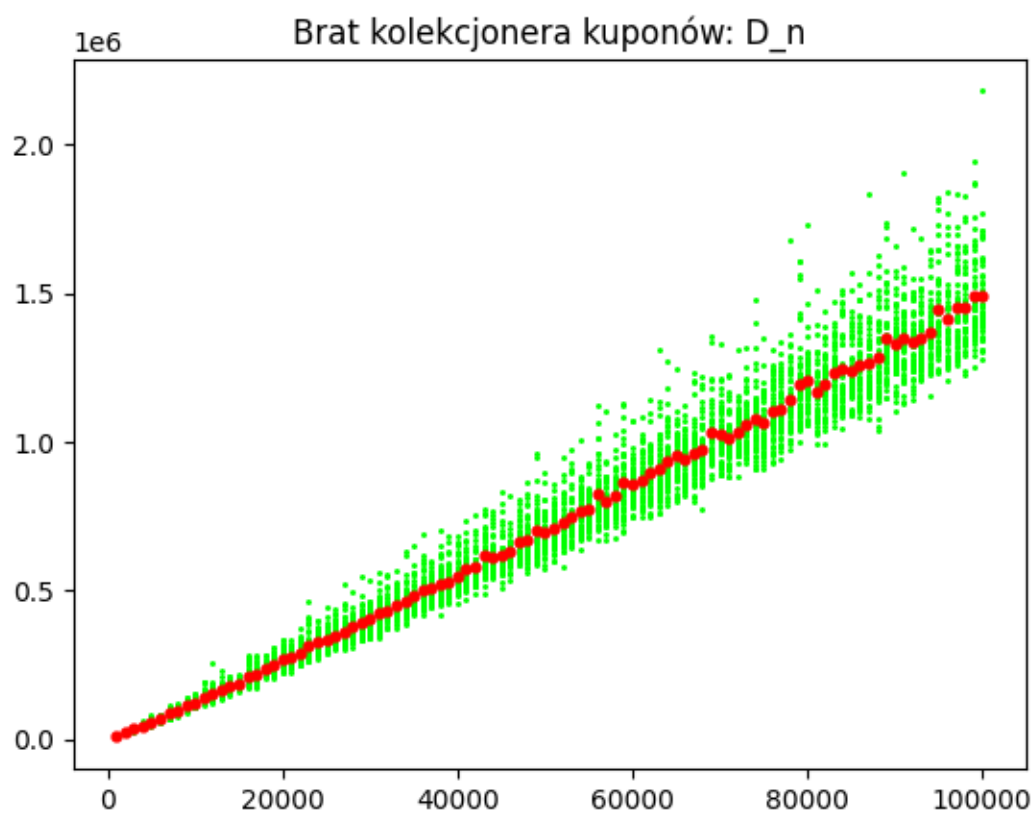


## Sprawozdanie do zadania 2

Program przeprowadził po  $k = 50$  niezależnych prób dla każdego  $n$  z przedziału  $\{1000, 2000, \dots, 100000\}$ , generując losowe liczby z przedziału  $\{1, 2, \dots, n\}$  odpowiadające "wrzuceniu kuli do  $n$ -tej urny" dopóki w każdej urnie nie znajdowały się przynajmniej 2 kule. Poniżej znajdują się wykresy obrazujące uzyskane wyniki dla różnych wielkości (czerwone kropki oznaczają wyniki średnie):





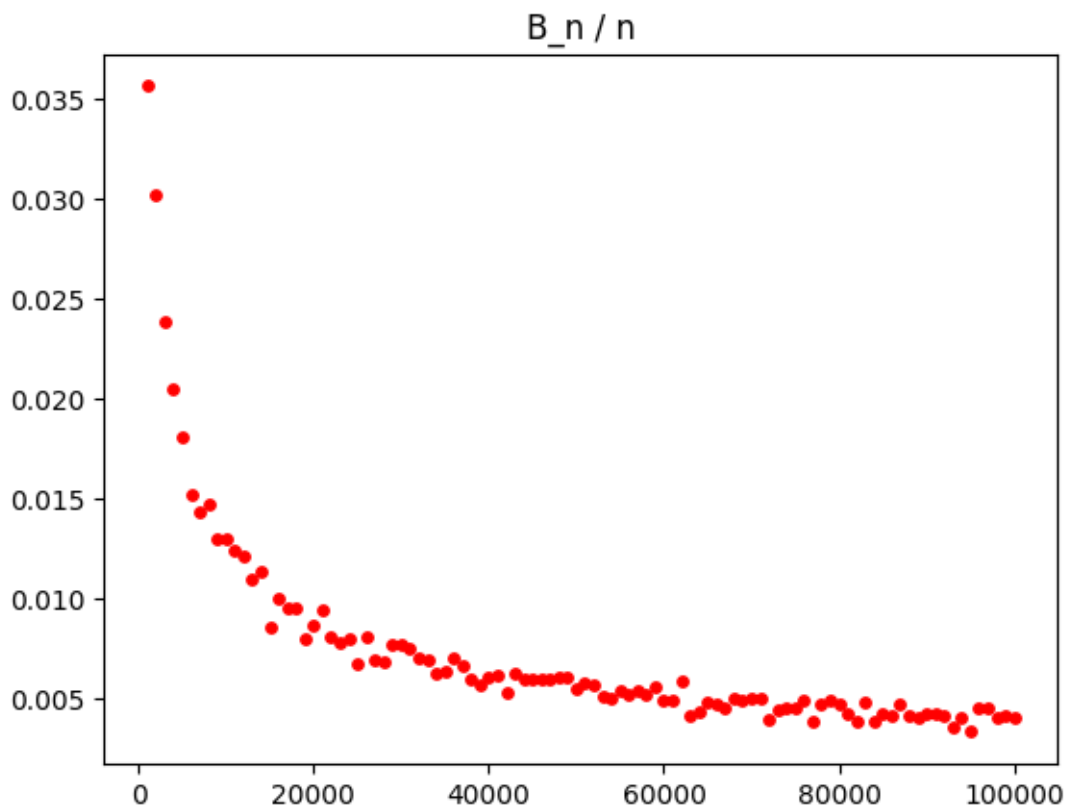


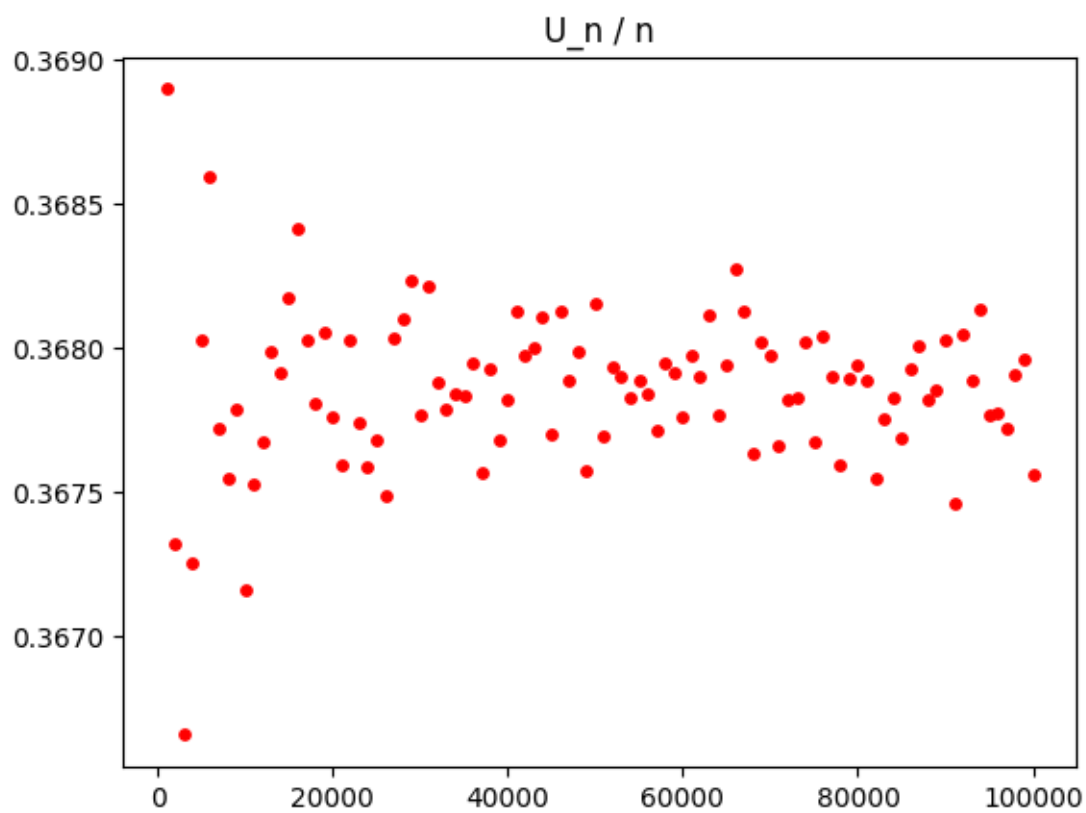
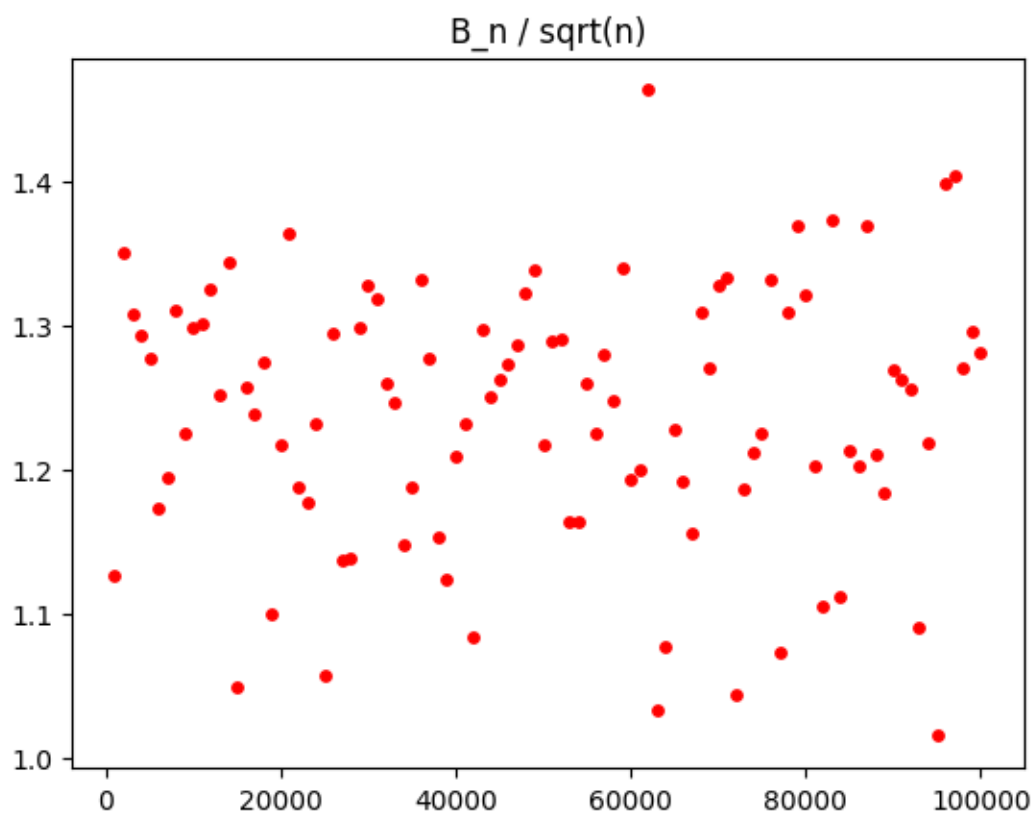
Mierzone wielkości to:

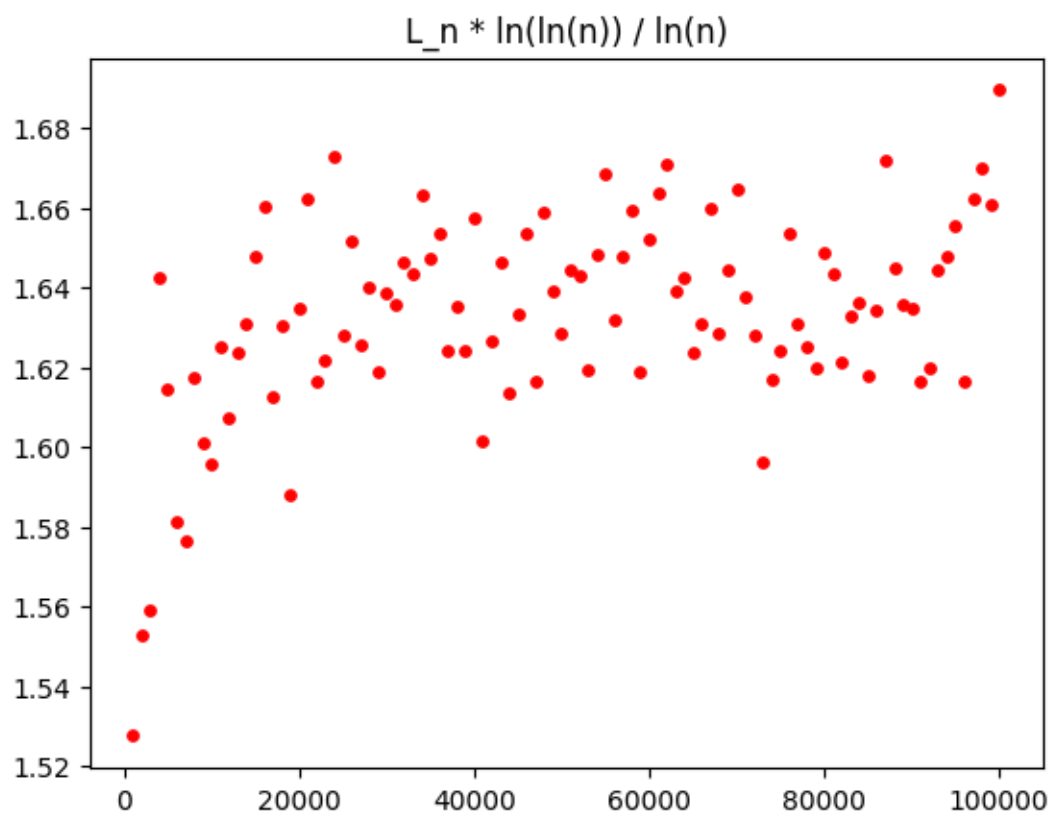
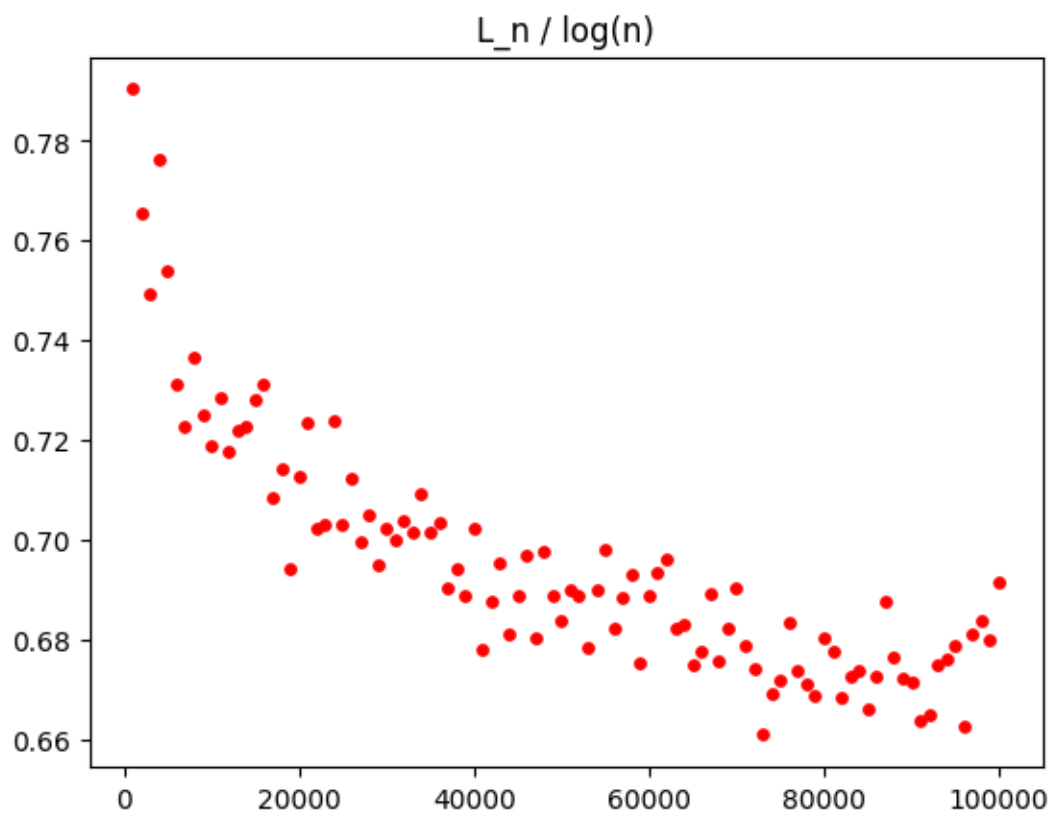
- $B_n$  – moment pierwszej kolizji, równy  $k$  jeśli  $k$ -ta z wrzucanych kul jest pierwszą, która trafiła do niepustej urny (paradoks urodzinowy)
- $U_n$  – liczba pustych urn po wrzuceniu  $n$  kul
- $L_n$  – maksymalna liczba kul w urnie po wrzuceniu  $n$  kul
- $C_n$  – minimalna liczba rzutów, po której w każdej z urn jest co najmniej jedna kula
- $D_n$  – minimalna liczba rzutów, po której w każdej z urn są co najmniej dwie kule
- $D_n - C_n$  – liczba rzutów od momentu  $C_n$  potrzebna do tego, żeby w każdej urnie były co najmniej dwie kule

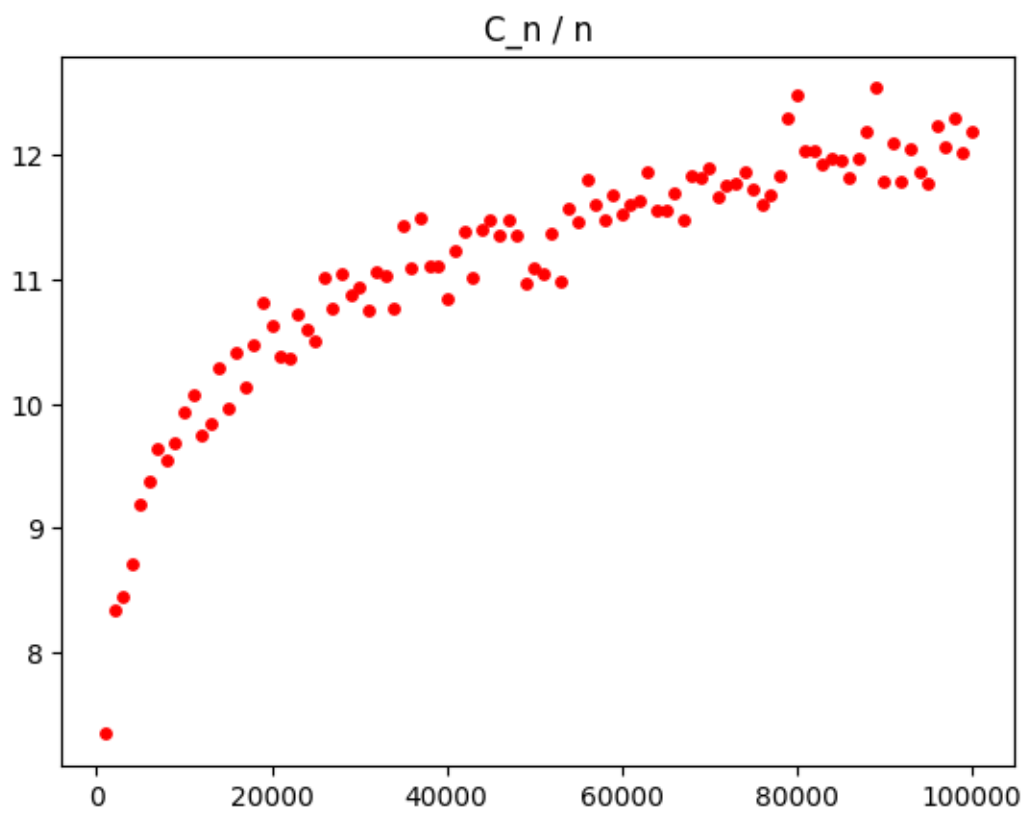
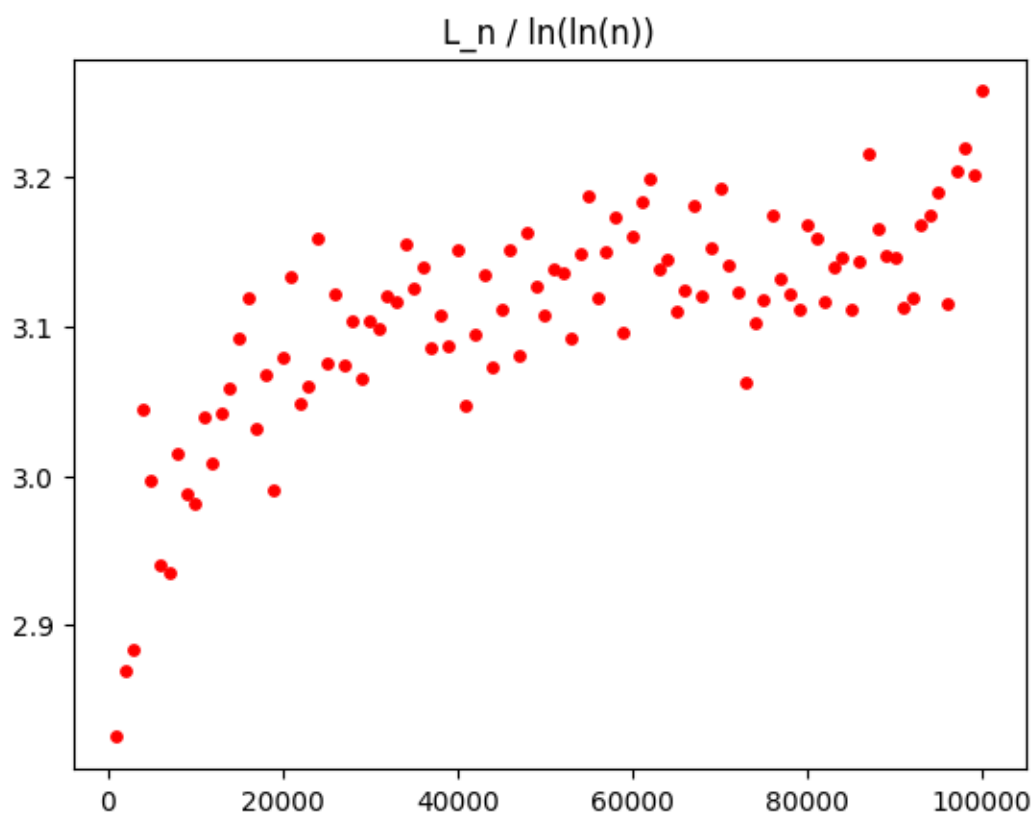
Jak widać z wykresów powyżej, wszystkie wartości średnie rosną wraz z wartością  $n$ , zgodnie z intuicją. Analizując wykres  $U_n$  dostrzegamy, że poszczególne punkty średnich wyznaczają z dużą dokładnością pewną prostą; ponadto, wartości dla przypadków poszczególnych mają tę samą własność. Dla wykresu  $L_n$  widać, że poszczególne wartości średnie są bardziej skoncentrowane dla większych wartości  $n$  – są one jednak słabo przybliżane przez wyniki poszczególnych eksperymentów (błąd względny pojedynczego wyniku może wynosić nawet około 25% nawet dla dużych  $n$ ). Dla wartości  $B_n$ ,  $C_n$ ,  $D_n$  oraz  $D_n - C_n$  obserwujemy, że precyzja estymacji zaczyna maleć wraz ze wzrostem  $n$  – średnia odległość poszczególnych średnich rośnie. To samo obserwujemy dla wyników poszczególnych doświadczeń – odchylenia od średnich wzrastają razem z wartością  $n$ .

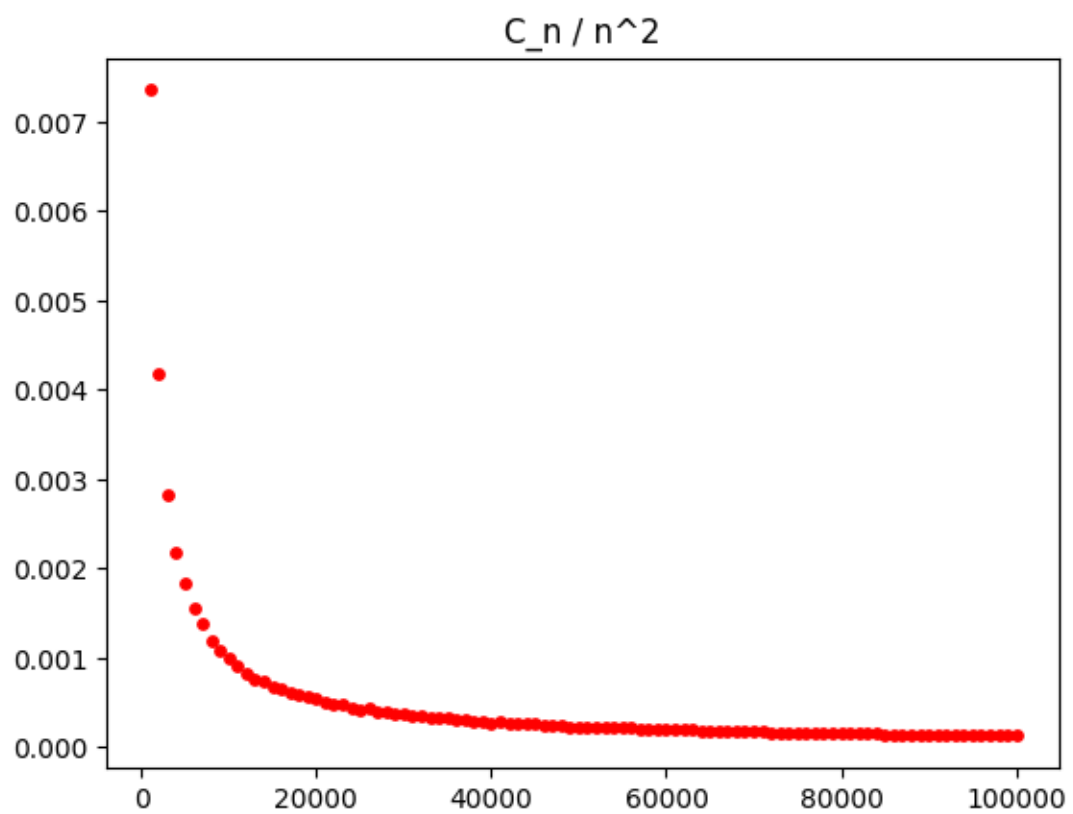
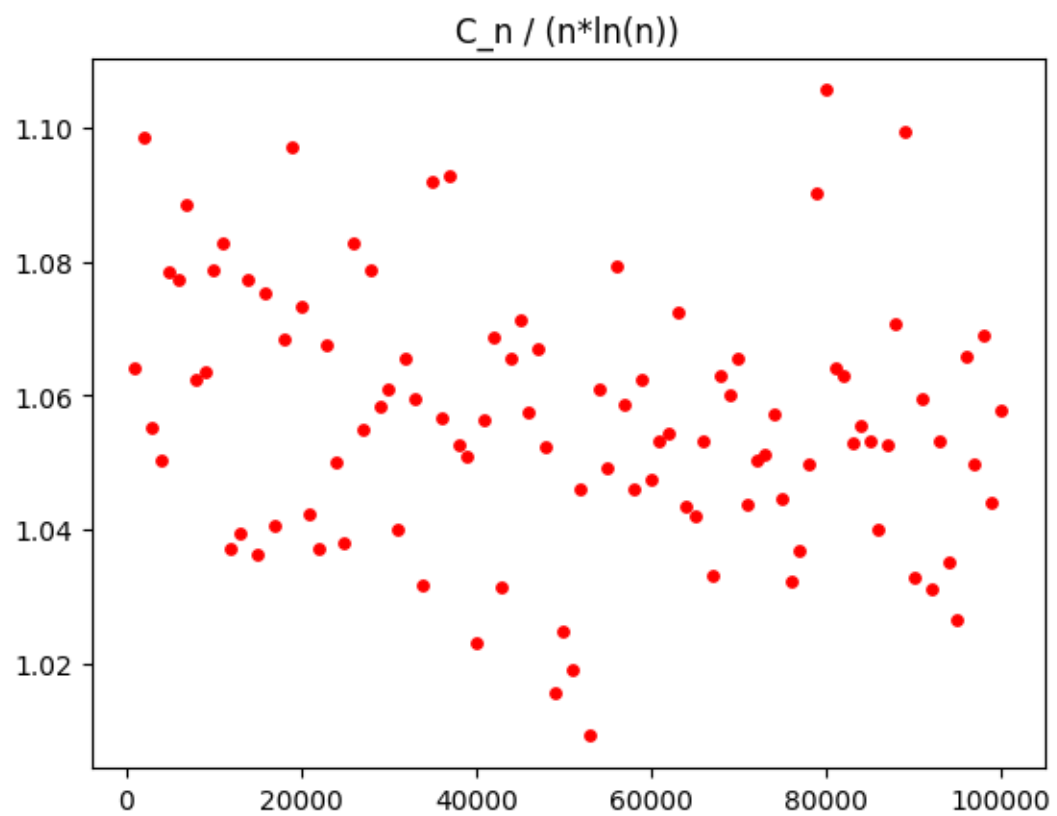
Poniższe wykresy ilustrują ilorazy średnich wartości poszczególnych wartości z wybranymi funkcjami zmiennej  $n$ :



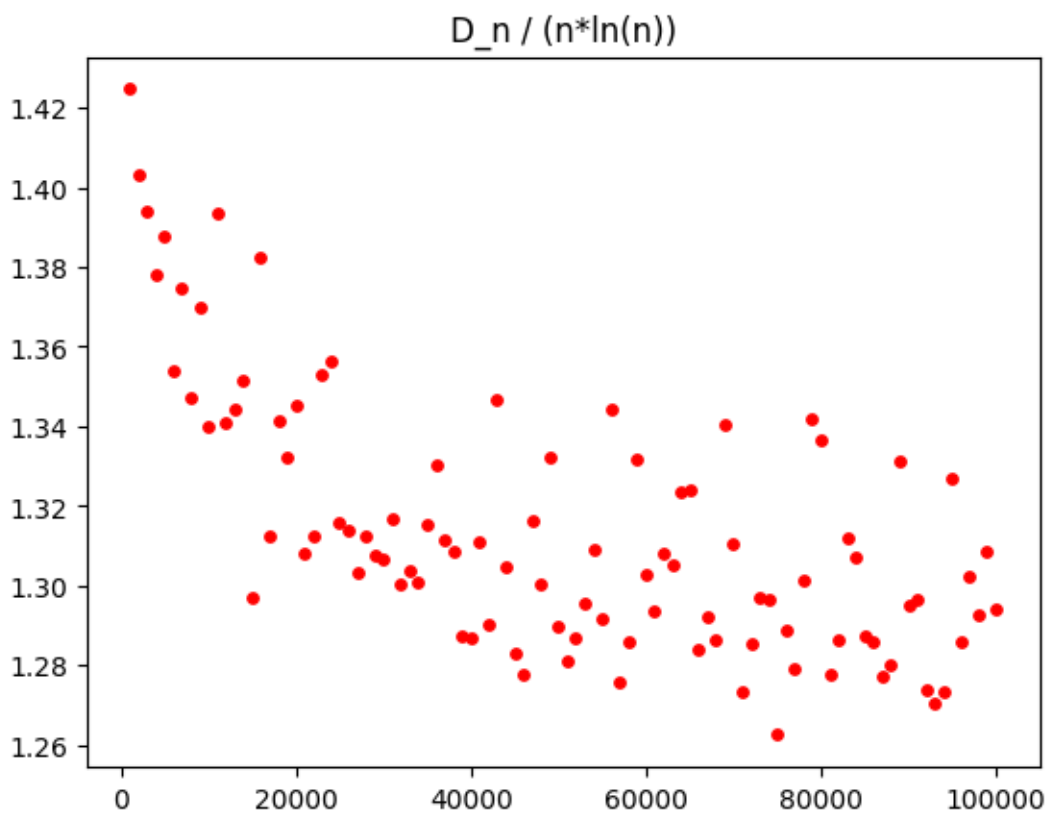
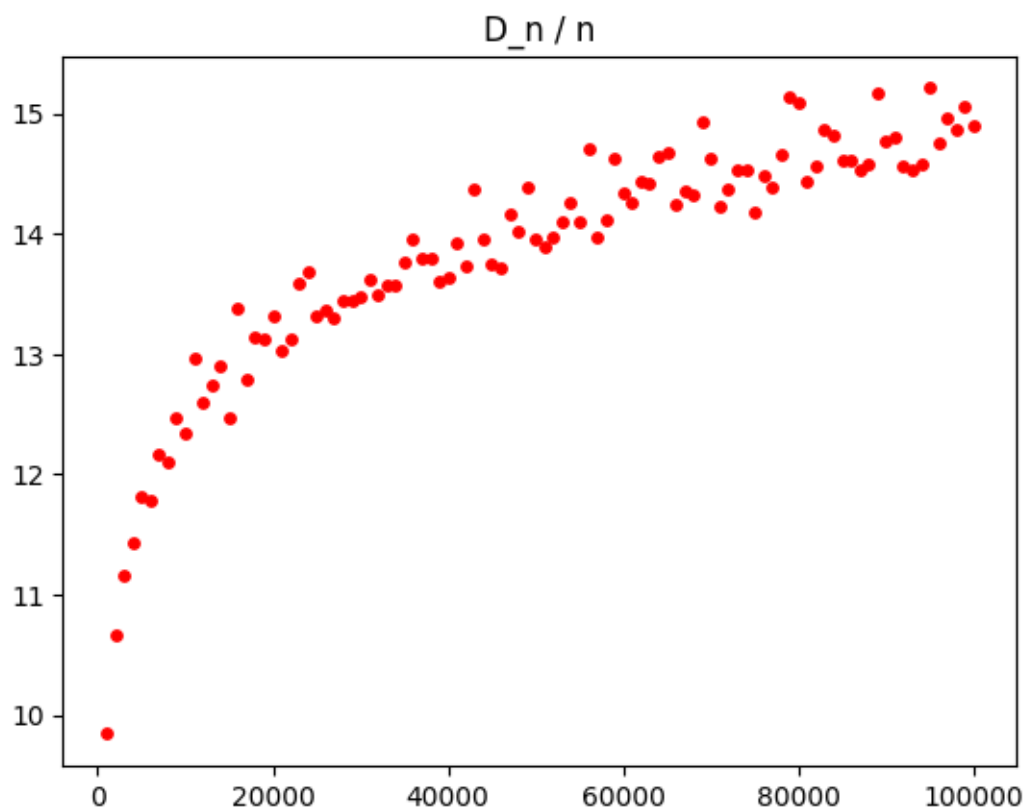


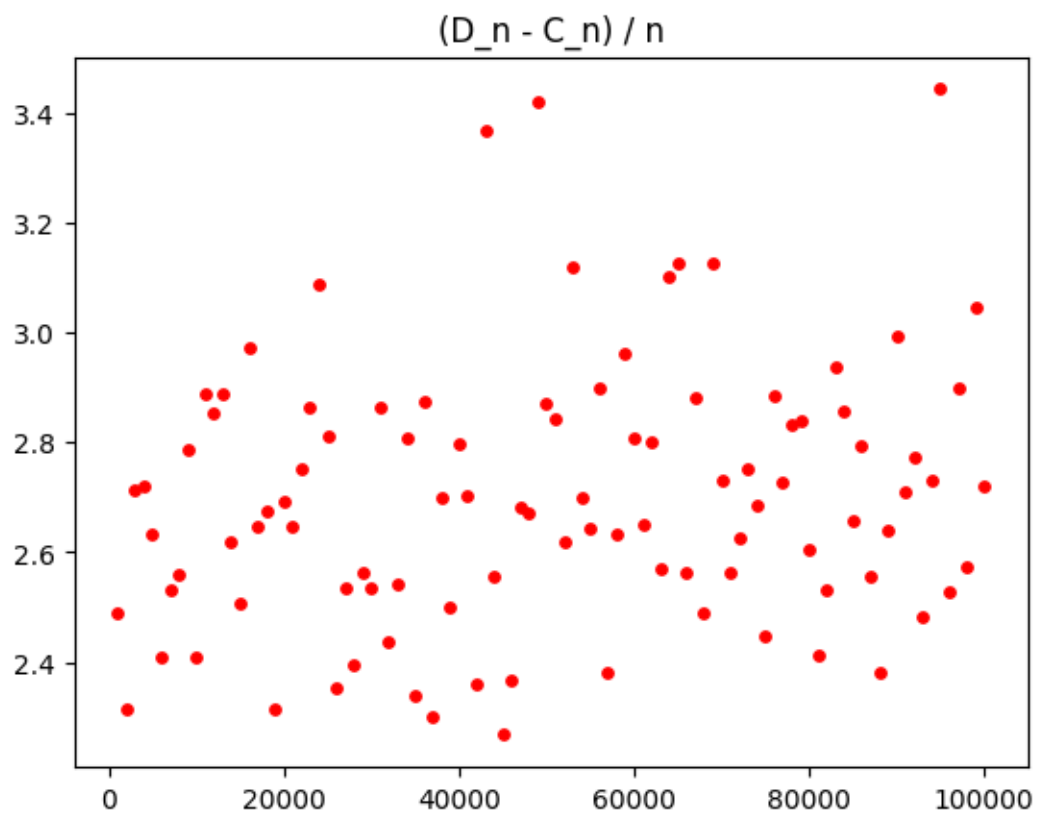
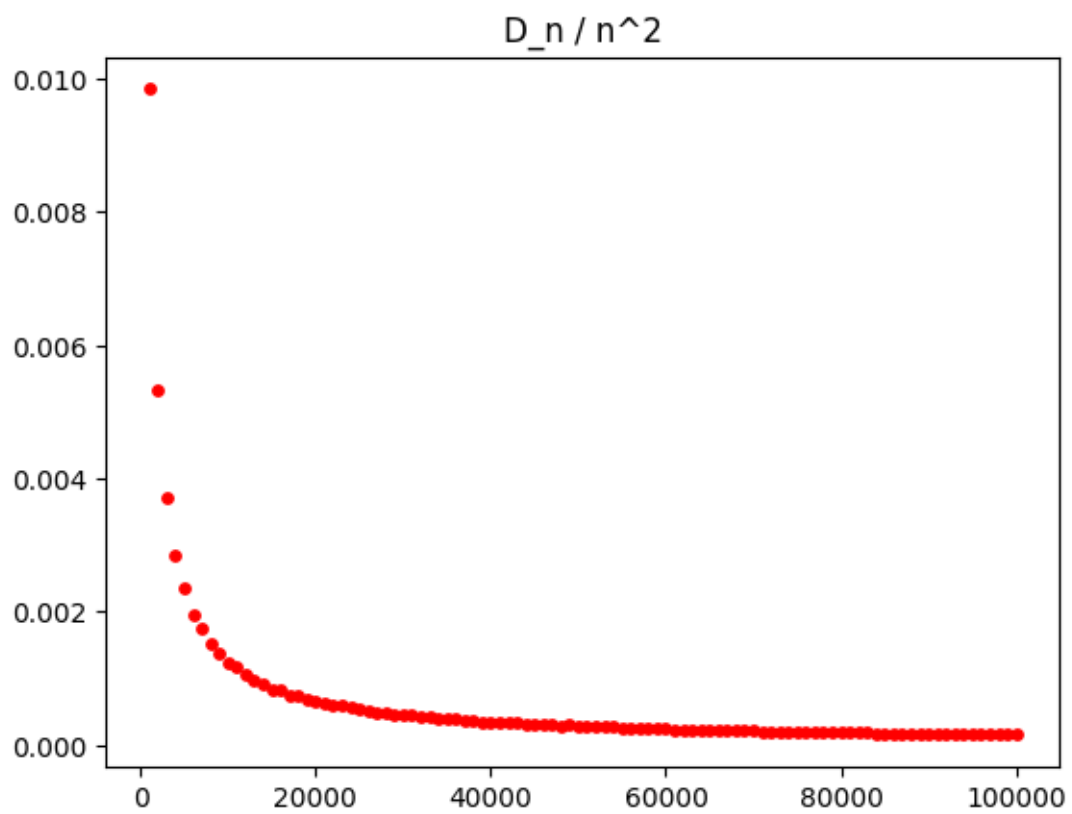


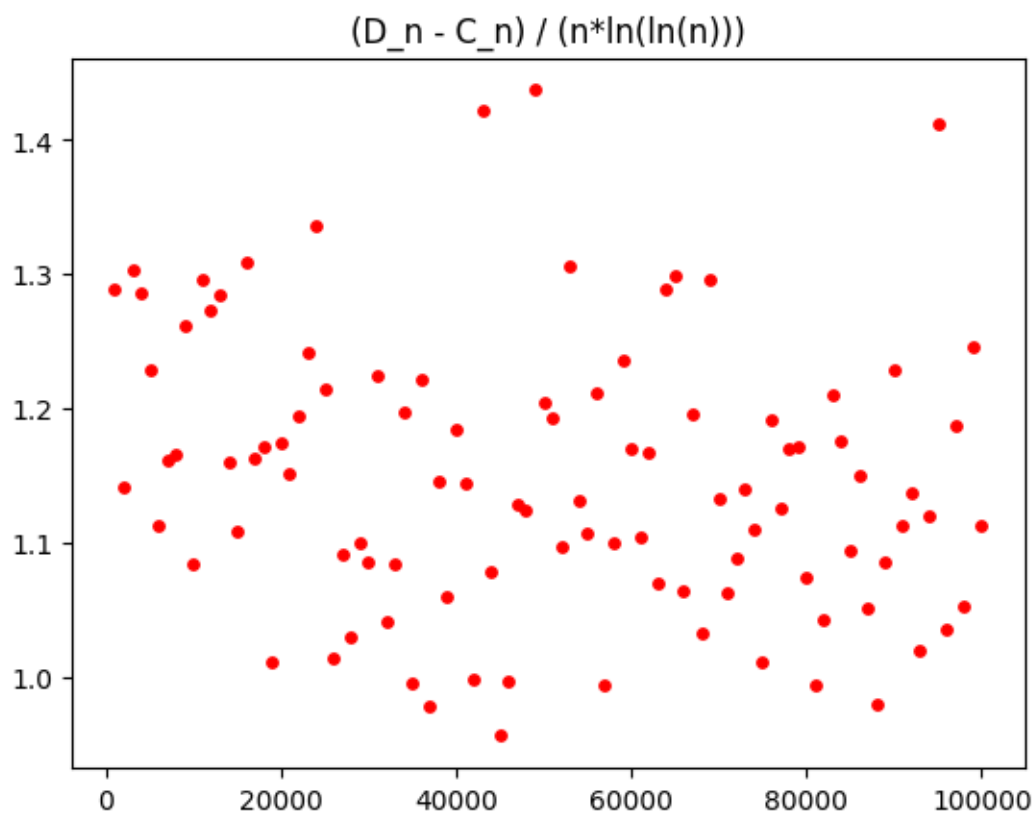
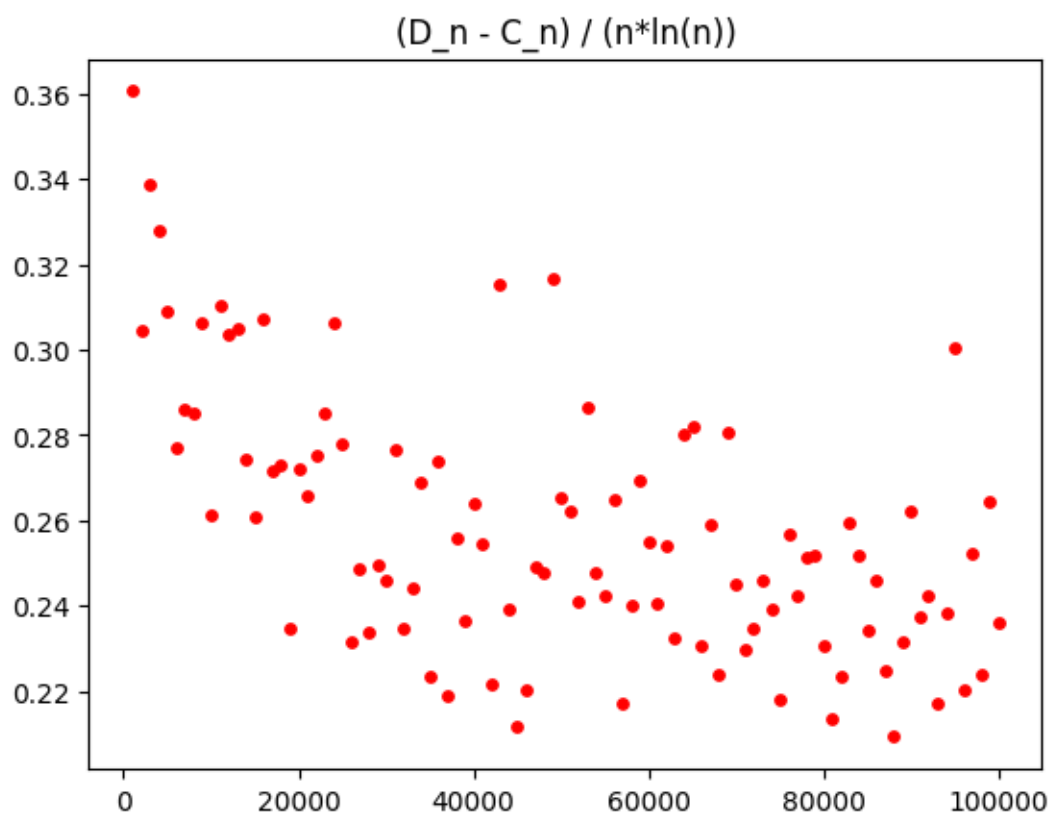












Na podstawie powyższych wykresów możemy postawić hipotezy dotyczące asymptotyki wyznaczonych wartości:

- $B_n = O(\sqrt{n})$
- $U_n = O(n)$
- $L_n = O(\ln(n) / \ln(\ln(n)))$
- $C_n = O(n \ln(n))$
- $D_n = O(n \ln(n))$
- $D_n - C_n = O(n \ln(\ln(n)))$

Nazwy "birthday paradox" oraz coupon collector's problem" można uzasadnić jak następuje:

- "birthday paradox" – wybieramy jednostajnie losowo po jednym człowieku z populacji na świecie. Czekamy na moment pierwszej kolizji, czyli pierwszy moment kiedy dwóch ludzi ma tę samą datę urodzenia. wynik jest paradoksalny z powodu tego, że uzyskane wyniki są nieintuicyjnie małe.
- "coupon collector's problem" – pewien człowiek zbiera kupony. Szukamy momentu, po którym będzie miał kupony każdego typu. Uzyskane wartość jest odpowiedzią na postawiony problem.

Paradoks urodzinowy jest ważny ze względu na funkcje haszujące, ponieważ zależy nam na dostępie do shaszowanych danych ze złożonością  $O(1)$ . Jeśli istnieją powtórzenia w wynikach haszy, to złożoność dostępu do danych rośnie. W kontekście kryptograficznych funkcji haszujących, musimy mieć pewność, że hasze haseł przechowywane w bazie danych są unikatowe – jeśli dwa różne hasła mają ten sam hasz, to dany program łamiący hasła musiałby średnio spędzić 2x krócej łamiąc hasło użytkownika, lub można byłoby się zalogować jako dany użytkownik korzystając z niepoprawnego hasła!